

[Open in app](#)

Parveen Khurana

124 Followers

[About](#)[Following](#)

Mathematics behind the parameters update rule:

 **Parveen Khurana** Jan 3, 2020 · 8 min read

This article covers the content discussed in the Sigmoid Neuron module of the [Deep Learning course](#) and all the images are taken from the same module.

In this article, we discuss the mathematics behind the parameters update rule.

Our goal is to find an algorithm which at any timestamp, tells us how to change the value of w such that the loss that we compute at the new value is less than the loss that we have at the current value.

$$w \rightarrow w + \eta \Delta w$$

$$\mathcal{L}(w) > \mathcal{L}(w + \eta \Delta w)$$

And if we keep doing this at every step, the loss is bound to decrease no matter where we start from and eventually reach its minimum value.

And Taylor series tells us that if we have a function and if we know its value at a certain point (x in the below case), then its value at a new point which is very close to x can be given by the below expression

[Open in app](#)


And we can see that the Taylor series relates the function value at a new point ($x + \delta x$) with the function value at the current point (x)

$$\underline{f(x + \Delta x)} = \underline{f(x)} + \left[\underline{f'(x)\Delta x} + \underline{\frac{1}{2!}f''(x)\Delta x^2} + \underline{\frac{1}{3!}f'''(x)\Delta x^3} + \dots \right]$$

In fact, the value at the new point is equal to the value at the current point plus some additional terms all of which depends on δx

Now if this δx is such that the quantity which is getting added to $f(x)$ (in brackets in the below image) is actually negative, then we can sure that the function value at a new point is less than the function value at the current point.

$$\underline{f(x + \Delta x)} = \underline{f(x)} + \left[\underline{f'(x)\Delta x} + \underline{\frac{1}{2!}f''(x)\Delta x^2} + \underline{\frac{1}{3!}f'''(x)\Delta x^3} + \dots \right]$$

~ve

So, we need to find δx in such a way so that the quantity in brackets in the above image is negative. The more negative the better because the loss would decrease by more.

$$f(x + \Delta x) = f(x) + \boxed{f'(x)}\Delta x + \frac{1}{2!}\boxed{f''(x)}\Delta x^2 + \frac{1}{3!}\boxed{f'''(x)}\Delta x^3 + \dots$$

The quantity in blue in the above image is the first-order derivative of x .

[Open in app](#)


The quantity in yellow in the above image is the third-order derivative of x .

If we have $f(x)$ as x^3 , then all of these quantities would be:

$$\begin{aligned} f(x) &= x^3 \\ f'(x) &= 3x^2 \\ f''(x) &= 6x \\ f'''(x) &= 6 \end{aligned}$$

Now we can write the Taylor series for the Loss function as:

$$\mathcal{L}(w + \Delta w) = \mathcal{L}(w) + \mathcal{L}'(w)\Delta w + \frac{1}{2!}\mathcal{L}''(w)\Delta w^2 + \frac{1}{3!}\mathcal{L}'''(w)\Delta w^3 + \dots$$

The idea is to find the value of Δw in such a way that the quantity in the brackets in the below image is negative, then we know that the new loss value would be smaller than the old loss

$$\mathcal{L}(w + \Delta w) = \mathcal{L}(w) + \left[\mathcal{L}'(w)\Delta w + \frac{1}{2!}\mathcal{L}''(w)\Delta w^2 + \frac{1}{3!}\mathcal{L}'''(w)\Delta w^3 \right] + \dots$$

[Open in app](#)


And since loss depends on b as well, so we want the new loss value to be less than the current loss value for new value of b as well

$$w \rightarrow w + \eta \Delta w$$

$$b \rightarrow b + \eta \Delta b$$

$$\mathcal{L}(w) > \mathcal{L}(w + \eta \Delta w)$$

$$\mathcal{L}(b) > \mathcal{L}(b + \eta \Delta b)$$

$$\mathcal{L}(w, b) > \mathcal{L}(w + \eta \Delta w, b + \eta \Delta b)$$

$$\mathcal{L}(\theta) > \mathcal{L}(\theta + \eta \Delta \theta) \quad \because \theta = [w, b]$$

If we change w or b , the predicted output would change.

If the predicted output changes, then the difference between the predicted output and the true output changes and if that is going to change, the loss value is going to change.

A handwritten diagram illustrating the vector update rule. It shows a vector θ (represented as a circle with a dot) being added to a scaled vector η multiplied by a column vector $\begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix}$. The entire expression is written as $\theta + \eta \begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix}$.

The Taylor series for a vector looks like:

Open in app



$$\Delta \theta = u$$

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \left[\eta * \underbrace{u^T \nabla_{\theta} \mathcal{L}(\theta)}_{\Delta \theta = u} + \frac{\eta^2}{2!} * \underbrace{u^T \nabla^2 \mathcal{L}(\theta) u}_{\Delta \theta = u} + \frac{\eta^3}{3!} * \dots + \dots \right]$$

The quantity in brackets in the above image depends on the change in the parameter value. So, we need to find this change vector (**$\Delta \theta$ or u**) such that the quantity in the brackets turns out to be negative and if that happens, we would be sure that the loss would decrease.

Now we can get rid of some of the terms in brackets in the below equation:

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta) + \left[\frac{\eta^2}{2!} * u^T \nabla^2 \mathcal{L}(\theta) u + \frac{\eta^3}{3!} * \dots + \dots \right]$$

Eta is very small, in practice, it would be something around 0.0001 or of that order, if **eta** is small then the higher powers of **eta** are going to be even smaller, so even without knowing the exact value of the quantity in yellow in the below image, we can be very sure that this entire term is going to be very small and similarly as we go ahead in the equation where we have **eta**³ and **eta**⁴, those terms would even be smaller

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta) + \left[\frac{\eta^2}{2!} * \underbrace{u^T \nabla^2 \mathcal{L}(\theta) u}_{\text{yellow}} + \frac{\eta^3}{3!} * \dots + \dots \right]$$

[Open in app](#)

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta) + \left[\frac{\eta^2}{2!} * u^T \nabla^2 \mathcal{L}(\theta) u + \frac{\eta^3}{3!} * \dots + \dots \right]$$

And we can write the equation as an approximation as the below:

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta) + \frac{\eta^2}{2!} * u^T \nabla^2 \mathcal{L}(\theta) u + \frac{\eta^3}{3!} * \dots + \dots$$

$$\mathcal{L}(\theta + \eta u) \approx \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta)$$

And the quantity in yellow in the above image is the first-order partial derivative of loss with respect to theta and the way to compute the partial derivative is that we assume the other variable to act as a constant for example: if the loss depends on w and b and we are computing the partial derivative with respect to w, then we can treat b as a constant.

If our function is below:

$$f(w, b) = w^3 + b^2$$

then we can compute the partial derivative with respect to 'w' as

$$\frac{\partial f(w, b)}{\partial w} = w^2$$

[Open in app](#)

and this would result in the below assuming **b** to be constant as we are computing partial derivative:

$$\frac{\partial f}{\partial w} \quad 3w^2$$

And similarly, the partial derivative with respect to **b** would be:

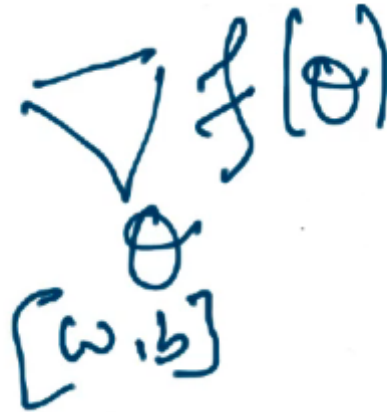
$$\frac{\partial f}{\partial b} \quad 2b$$

And if we put the above two partial derivatives(with respect to **w** and **b**) in a vector, we get a gradient:

$$\begin{bmatrix} \frac{\partial f}{\partial w} \\ \frac{\partial f}{\partial b} \end{bmatrix} \begin{bmatrix} 3w^2 \\ 2b \end{bmatrix}$$

[Open in app](#)


And for the above case, we can write as:



which means it is the gradient of the function $f(\theta)$ with respect to θ and θ actually is just a vector of w and b .

So, going back to our original equation, we have:

$$\underbrace{\mathcal{L}(\theta + \eta u)}_{\mathbb{R}} \approx \underbrace{\mathcal{L}(\theta)}_{\mathbb{R}} + \underbrace{\eta}_{\mathbb{R}} * \underbrace{u^T \nabla_{\theta} \mathcal{L}(\theta)}_{\mathbb{R}}$$

$\begin{bmatrix} \Delta w & \Delta b \end{bmatrix} \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial w} \\ \frac{\partial \mathcal{L}}{\partial b} \end{bmatrix}$

The quantity on the left-hand side is the loss value at the new point and is going to a Real no., the current loss value is also a real no., η is also a small real number, and

[Open in app](#)


$$\underbrace{\mathcal{L}(\theta + \eta u)}_{\mathbb{R}} \approx \underbrace{\mathcal{L}(\theta)}_{\mathbb{R}} + \underbrace{\eta}_{\mathbb{R}} * \underbrace{u^T \nabla_{\theta} \mathcal{L}(\theta)}_{\mathbb{R}}$$

We can re-write the above equation as:

$$\mathcal{L}(\theta + \eta u) = \mathcal{L}(\theta) + \eta * u^T \nabla_{\theta} \mathcal{L}(\theta)$$

$$\mathcal{L}(\theta + \eta u) - \mathcal{L}(\theta) = \eta * u^T \nabla_{\theta} \mathcal{L}(\theta)$$

Note that the move ηu would be favorable only if,

$$\mathcal{L}(\theta + \eta u) - \mathcal{L}(\theta) < 0 \quad [\text{i.e. if the new loss is less than the previous loss}]$$

This implies,

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0$$

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0$$

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0$$

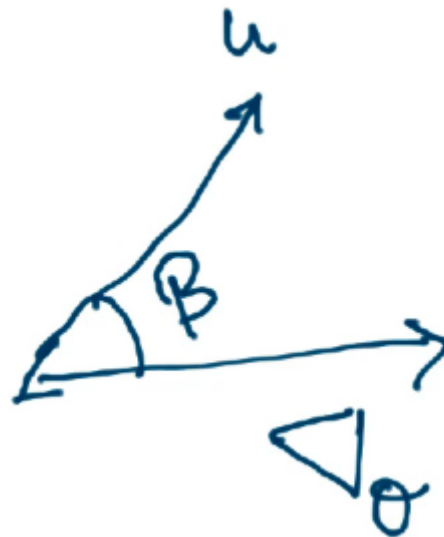
The quantity in the above image is a dot product between two vectors.

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0 \rightarrow \mathbb{R}^2$$

[Open in app](#)

 \mathbb{R}^2

$$\begin{bmatrix} \Delta w \\ \Delta b \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{bmatrix}$$



$$\cos \beta = \frac{u^T \nabla_{\theta} L(\theta)}{\|u\| \|\nabla_{\theta} L(\theta)\|}$$

[Open in app](#)


We want the below quantity to be less than 0 and since it is the dot product of two vectors, we want the angle between these two vectors to be greater than 90 but less than equal to 180:

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0$$

Okay, so we have,

$$u^T \nabla_{\theta} \mathcal{L}(\theta) < 0$$

Let β be the angle between u and $\nabla_{\theta} \mathcal{L}(\theta)$, then we know that,

$$-1 \leq \cos(\beta) = \frac{u^T \nabla_{\theta} \mathcal{L}(\theta)}{\|u\| * \|\nabla_{\theta} \mathcal{L}(\theta)\|} \leq 1$$

multiply throughout by $k = \|u\| * \|\nabla_{\theta} \mathcal{L}(\theta)\|$

$$-k \leq u^T \nabla_{\theta} \mathcal{L}(\theta) \leq k$$

Thus, $\mathcal{L}(\theta + \eta u) - \mathcal{L}(\theta) = u^T \nabla_{\theta} \mathcal{L}(\theta) = k * \cos \beta$ will be most negative when $\cos(\beta) = -1$ i.e., when β is 180°

Gradient Descent Rule,

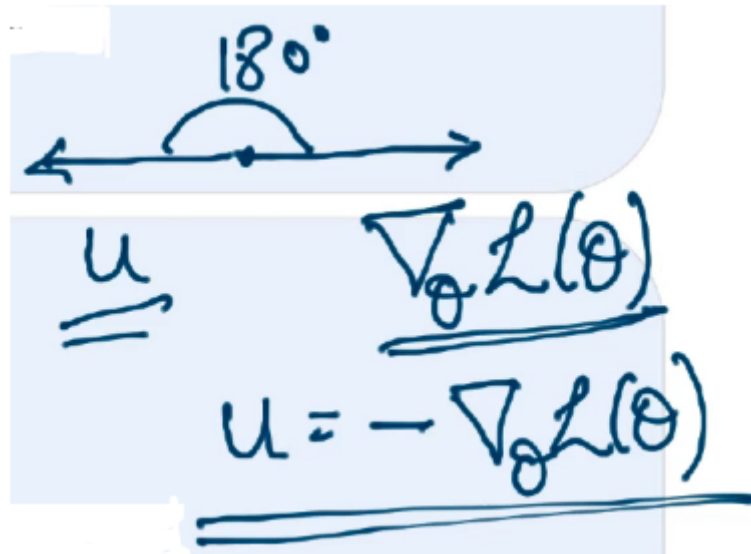
- The direction u that we intend to move in should be at 180° w.r.t. the gradient.
- In other words, move in a direction opposite to the gradient.

Parameter Update Rule

$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

$$\text{where } \Delta w_t = \frac{\partial \mathcal{L}(w, b)}{\partial w} \text{ at } w=w_t, b=b_t, \Delta b_t = \frac{\partial \mathcal{L}(w, b)}{\partial b} \text{ at } w=w_t, b=b_t$$

[Open in app](#)


Computing Partial Derivatives:

The general recipe that we have is:

Initialise w, b

Iterate over data:

compute \hat{y}

compute $\mathcal{L}(w, b)$

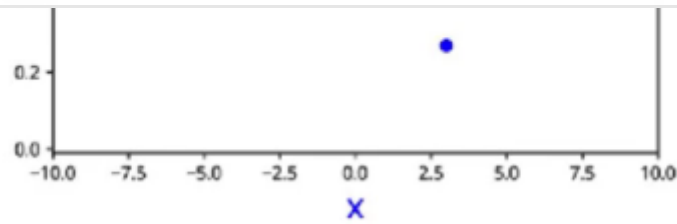
$$w_{t+1} = w_t - \eta \Delta w_t$$

$$b_{t+1} = b_t - \eta \Delta b_t$$

till satisfied

So, we have 5 data points as the input, we have chosen the model to be Sigmoid function



[Open in app](#)


And we compute the loss value as:

$$\mathcal{L} = \frac{1}{5} \sum_{i=1}^{i=5} \underbrace{(f(x_i) - y_i)}_{\text{residual}}^2$$

We can compute the δw as:

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i) \right]$$

$$\Delta w = \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{5} \sum_{i=1}^{i=5} \frac{\partial}{\partial w} (f(x_i) - y_i)$$

Consider the square in this equation as the loss function we are using is squared error loss.

The derivative of a sum of quantities is equal to the sum of the derivative of individual quantities.

And now out of the 5 terms in the derivative, let's consider one term and we will compute its partial derivative with respect to w :

$$\Delta w = \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{5} \sum_{i=1}^{i=5} \frac{\partial}{\partial w} (f(x_i) - y_i)$$

Let's consider only one term in this sum

[Open in app](#)


Consider the square in this equation as the loss function we are using is squared error loss.

Considering one term we have:

$$\nabla w = \frac{\partial}{\partial w} \left[\frac{1}{2} * (f(x) - y)^2 \right]$$

We have taken 1/2 in the above equation just for the sake of convenience.

We have $f(x)$ as a function of 'w' as $f(x)$ equals:

$$\frac{1}{1 + e^{-wx + b}}$$

So, using the chain rule, we can compute the partial derivative with respect to 'w' as:

$$\begin{aligned} \nabla w &= \frac{\partial}{\partial w} \left[\frac{1}{2} * (f(x) - y)^2 \right] \\ &= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y)] \end{aligned}$$

Now y in the above equation does not depend on w , so its partial derivative with respect to w would be 0 and we can write the above equation as:

$$\nabla w = \frac{\partial}{\partial w} \left[\frac{1}{2} * (f(x) - y)^2 \right]$$

[Open in app](#)


$$\begin{aligned}
 &= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w}(f(x) - y)] \\
 &= (f(x) - y) * \frac{\partial}{\partial w}(f(x))
 \end{aligned}$$

And now we can plug in the value of $f(x)$ in the above equation, so we have:

$$\begin{aligned}
 \nabla w &= \frac{\partial}{\partial w} [\frac{1}{2} * (f(x) - y)^2] \\
 &= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w}(f(x) - y)] \\
 &= (f(x) - y) * \frac{\partial}{\partial w}(f(x)) \\
 &= (f(x) - y) * \frac{\partial}{\partial w} \left(\frac{1}{1 + e^{-(wx+b)}} \right)
 \end{aligned}$$

And we can compute the partial derivative of $f(x)$ with respect to w as:

$$\begin{aligned}
 &\frac{\partial}{\partial w} \left(\frac{1}{1 + e^{-(wx+b)}} \right) \\
 &= \frac{-1}{(1 + e^{-(wx+b)})^2} \frac{\partial}{\partial w} (e^{-(wx+b)}) \\
 &= \frac{-1}{(1 + e^{-(wx+b)})^2} * (e^{-(wx+b)}) \frac{\partial}{\partial w} (-(wx + b)) \\
 &= \frac{-1}{(1 + e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1 + e^{-(wx+b)})} * (-x) \\
 &= \frac{1}{(1 + e^{-(wx+b)})} * \frac{e^{-(wx+b)}}{(1 + e^{-(wx+b)})} * (x)
 \end{aligned}$$

[Open in app](#)


And our overall partial derivative value looks like:

$$\begin{aligned}
 \nabla w &= \frac{\partial}{\partial w} \left[\frac{1}{2} * (f(x) - y)^2 \right] \\
 &= \frac{1}{2} * [2 * (f(x) - y) * \frac{\partial}{\partial w} (f(x) - y)] \\
 &= (f(x) - y) * \frac{\partial}{\partial w} (f(x)) \\
 &= (f(x) - y) * \frac{\partial}{\partial w} \left(\frac{1}{1 + e^{-(wx+b)}} \right) \\
 &= (f(x) - y) * f(x) * (1 - f(x)) * x
 \end{aligned}$$

$$\mathcal{L} = \frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i)$$

$$\frac{\partial \mathcal{L}}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{5} \sum_{i=1}^{i=5} (f(x_i) - y_i) \right]$$

$$\Delta w = \frac{\partial \mathcal{L}}{\partial w} = \frac{1}{5} \sum_{i=1}^{i=5} \frac{\partial}{\partial w} (f(x_i) - y_i)$$

Let's consider only one term in this sum,

$$\Delta w = (f(x) - y) * f(x) * (1 - f(x)) * x$$

For 5 points,

[Open in app](#)

So, this is how we compute the partial derivative of the loss function with respect to the parameter ' \mathbf{w} '. In the same manner, we can compute the partial derivative of the loss function with respect to the parameter ' \mathbf{b} '.

[Machine Learning](#)[Artificial Intelligence](#)[Deep Learning](#)[Sigmoid](#)[Artificial Neuron](#)[About](#) [Write](#) [Help](#) [Legal](#)

Get the Medium app

